Natural Language Processing

# Part-Of-Speech Tagging

LESSON 18

prof. Antonino Staiano

M.Sc. In ''Machine Learning e Big Data'' - University Parthenope of Naples
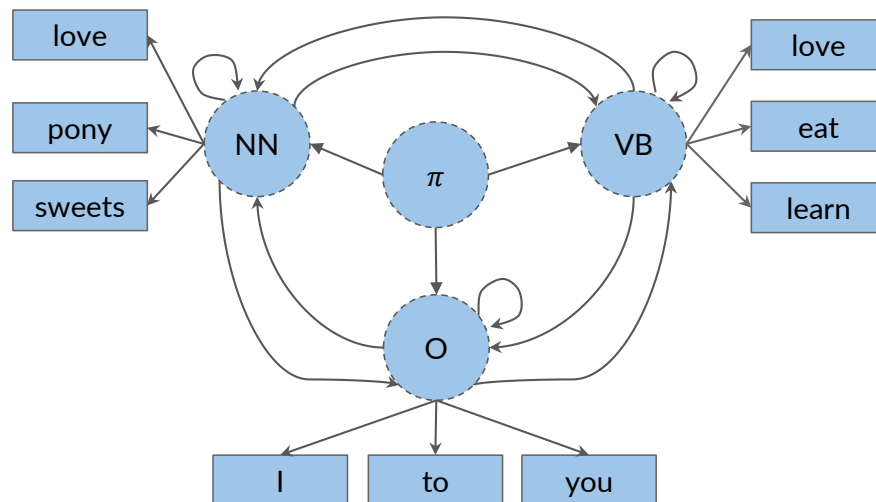
HMM for POS Tagging

# The Viterbi Algorithm

# HMM tagging as decoding

- The task of determining the sequence of the hidden variables corresponding to the sequence of observations is called decoding

- Decoding
  - *Given as input an HMM with A and B matrices, and a sequence of observations $O=o_1,o_2, \ldots, o_T$, find the most probable sequence of states $Q=q_1q_2q_3\ldots q_T$*
  - For POS tagging, the goal of HMM decoding is to choose the tag sequence $t_1,\ldots,t_n$ that is most probable given the observation sequence of n words $w_1,\ldots,w_n$

- The decoding algorithm for HMMs is the Viterbi algorithm
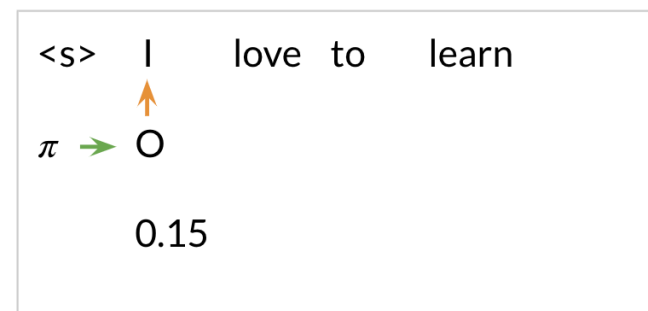
# Viterbi Algorithm: The big picture
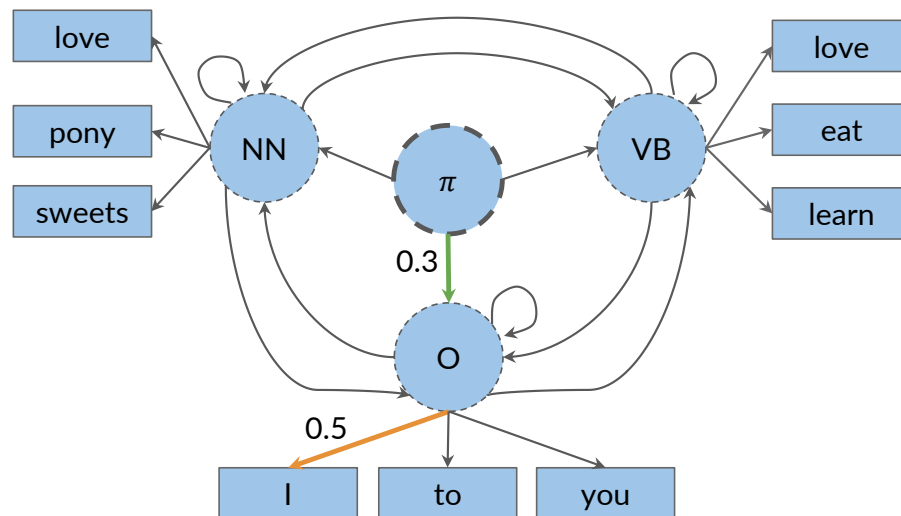
- Find the sequence of hidden states or parts of speech tags that have the highest probability for this sequence

# Viterbi Algorithm: The big picture

- Let's start from the initial state $\pi$, selecting the next most probable hidden state



- The joint probability for observing the word *I* and with a transition through the O state is 0.15 (0.3 x 0.5, i.e., transition prob x emission prob)

# Viterbi Algorithm: The big picture

- Now, two possibilities of having observed the word *love*

# Viterbi Algorithm: The big picture

# Viterbi Algorithm: The big picture

# Viterbi Algorithm: The big picture

# Viterbi Algorithm: The big picture

- The total probability is the product of all the probabilities for the single steps chosen



Probability for this sequence of hidden states: 0.0003

- The Viterbi algorithm computes several such paths at the same time for finding the most likely sequence of hidden states

# Viterbi algorithm: Steps

- Initialization step
- Forward pass
- Backward pass

$$C = \begin{array}{|c|c|c|c|c|} \hline & w_1 & w_2 & ... & w_K \\ \hline t_1 & & & & \\ \hline ... & & & & \\ \hline t_N & & & & \\ \hline \end{array}$$

$$D = \begin{array}{|c|c|c|c|c|} \hline & w_1 & w_2 & ... & w_K \\ \hline t_1 & & & & \\ \hline ... & & & & \\ \hline t_N & & & & \\ \hline \end{array}$$

- Auxiliary matrices
  - C holds the intermediate optimal probabilities
  - D holds the indices of the visited states
  - size NxK, N= number of POS tags, K = number of words in the given sequence

UNIVERSITÀ DEGLI STUDI DI NAPOLI
PARTHENOPE

# Initialization step: C matrix

- The first columns of matrices C and D are populated
  - C
    - The first column represents the probability of the transitions from the start state $\pi$ to the first tag_i and the word $w_1$



$$C = \quad {}^t_{..}C = \quad {}^t$$

| t | $w_1$ | $w_2$ | ... | $w_K$ |
|---|---|---|---|---|
| $t_1$ | $c_{1,1}$ | | | |
| ... | | | | |
| $t_N$ | $c_{N,1}$ | | | |

$$c_{i,1} = \boxed{\pi_i} * \boxed{b_{i,cindex(w_1)}}$$
$$= a_{1,i} * b_{i,cindex(w_1)}$$

- *cindex($w_i$)* returns the column index in the emission matrix B for a given word, $w_i$

# Initialization step: D matrix

- D stores the labels that represent the different states we're traversing when finding the most likely sequence of POS tags from the given sequence of words, from $w_1$ to $w_k$
  - The first column has all 0 entries as there are no preceding POS tags traversed



$$D = \begin{array}{c} t \\ \\ \cdots \\ \\ t \end{array}$$

$$D =$$

|  | $w_1$ | $w_2$ | ... | $w_K$ |
|---|---|---|---|---|
| $t_1$ | $d_{1,1}$ |  |  |  |
| ... |  |  |  |  |
| $t_N$ | $d_{N,1}$ |  |  |  |

$$d_{i,1} = 0 \qquad d_{i,1} = 0$$

# Forward pass

- C and D are populated column by column during the forward pass



| | $w_1$ | $w_2$ | ... | $w_K$ |
|---|---|---|---|---|
| $t_1$ | $c_{1,1}$ | $c_{1,2}$ | | $c_{1,K}$ |
| ... | | | | |
| $t_N$ | $c_{N,1}$ | $c_{N,2}$ | | $c_{N,K}$ |

$$D:$$
$$C =$$

$$c_{i,j} = \max_k c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

# Forward pass



| | $w_1$ | $w_2$ | ... | $w_K$ |
|---|---|---|---|---|
| $t_1$ | $c_{1,1}$ | $c_{1,2}$ | | $c_{1,K}$ |
| ... | | | | |
| $t_N$ | $c_{N,1}$ | $c_{N,2}$ | | $c_{N,K}$ |

$$C = \quad {}_t C =$$

$$c_{1,2} = \max_k c_{k,1} * a_{k,1} * \boxed{b_{1,cindex(w_2)}}$$

$$c_{1,2} = \max_k c_{k,1} * a_{k,1} * b_{1,cindex(w_2)}$$

# Forward pass



| t | $w_1$ | $w_2$ | ... | $w_K$ |
|---|---|---|---|---|
| $t_1$ | $c_{1,1}$ | $c_{1,2}$ | | $c_{1,K}$ |
| ... | | | | |
| $t_N$ | $c_{N,1}$ | $c_{N,2}$ | | $c_{N,K}$ |

$$C = $$
$$C = $$

$$c_{1,2} = \max_k c_{k,1} * a_{k,1} * b_{1,cindex(w_2)}$$
$$c_{1,2} = \max_k c_{k,1} * a_{k,1} * b_{1,cindex(w_2)}$$

# Forward pass



$$C = \begin{matrix} t \\ \\ \dots \\ \\ t \end{matrix}$$

| | $w_1$ | $w_2$ | ... | $w_K$ |
|---|---|---|---|---|
| $t_1$ | $c_{1,1}$ | $c_{1,2}$ | | $c_{1,K}$ |
| ... | | | | |
| $t_N$ | $c_{N,1}$ | $c_{N,2}$ | | $c_{N,K}$ |

$$c_{1,2} = \max_k c_{1,2} = \max_k \boxed{c_{k,1}} * \boxed{a_{k,1}} * \boxed{b_{1,cindex(w_2)}}$$

# Forward pass



|     | $w_1$ | $w_2$ | ... | $w_K$ |
|-----|-------|-------|-----|-------|
| $t_1$ | $d_{1,1}$ | $d_{1,2}$ |  | $d_{1,K}$ |
| ... |  |  |  |  |
| $t_N$ | $d_{N,1}$ | $d_{N,2}$ |  | $d_{N,K}$ |

$$D =$$

$$c_{i,j} = \max_k c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

$$d_{i,j} = \underset{k}{\operatorname{argmax}}\, c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

- In each $d_{i,j}$ the k which maximizes the entry $c_{i,j}$ is stored

# Backward pass

- The forward pass provided us the matrix C and D populated

$$C = \begin{array}{|c|c|c|c|c|} \hline & w_1 & w_2 & \cdots & w_K \\ \hline t_1 & c_{1,1} & c_{1,2} & & c_{1,K} \\ \hline \cdots & & & & \\ \hline t_N & c_{N,1} & c_{N,2} & & c_{N,K} \\ \hline \end{array} \qquad D = \begin{array}{|c|c|c|c|c|} \hline & w_1 & w_2 & \cdots & w_K \\ \hline t_1 & d_{1,1} & d_{1,2} & & d_{1,K} \\ \hline \cdots & & & & \\ \hline t_N & d_{N,1} & d_{N,2} & & d_{N,K} \\ \hline \end{array}$$

$$s = \operatorname*{argmax}_i c_{i,K}$$

# Backward pass

- First, calculate the index of the entry $c_{iK}$ in the last column of C

$$C = \begin{array}{c|c|c|c|c|} & w_1 & w_2 & \ldots & w_K \\ \hline t_1 & c_{1,1} & c_{1,2} & & c_{1,K} \\ \hline \ldots & & & & \\ \hline t_N & c_{N,1} & c_{N,2} & & c_{N,K} \\ \hline \end{array}$$

$$D = \begin{array}{c|c|c|c|c|} & w_1 & w_2 & \ldots & w_K \\ \hline t_1 & d_{1,1} & d_{1,2} & & d_{1,K} \\ \hline \ldots & & & & \\ \hline t_N & d_{N,1} & d_{N,2} & & d_{N,K} \\ \hline \end{array}$$

$$s = \operatorname*{argmax}_{i} c_{i,K}$$

# Backward pass

- Example

$$D = \begin{array}{c|c|c|c|c|c|}
 & \text{w}_1 & \text{w}_2 & \text{w}_3 & \text{w}_4 & \text{w}_5 \\
\hline
\text{t}_1 & 0 & 1 & 3 & 2 & 3 \\
\hline
\text{t}_2 & 0 & 2 & 4 & 1 & 3 \\
\hline
\text{t}_3 & 0 & 2 & 4 & 1 & 4 \\
\hline
\text{t}_4 & 0 & 4 & 4 & 3 & 1 \\
\hline
\end{array}$$

| \<s\> | w1 | w2 | w3 | w4 | w5 |

# Backward pass

$$C =$$

| | $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ |
|---|---|---|---|---|---|
| $t_1$ | 0.25 | 0.125 | 0.025 | 0.0125 | 0.01 |
| $t_2$ | 0.1 | 0.025 | 0.05 | 0.01 | 0.003 |
| $t_3$ | 0.3 | 0.05 | 0.025 | 0.02 | 0.0000 |
| $t_4$ | 0.2 | 0.1 | 0.000 | 0.0025 | 0.0003 |

$$s = \operatorname*{argmax}_{i} c_{i,K} = 1$$

# Backward pass

$$D = \begin{array}{c|ccccc} & w_1 & w_2 & w_3 & w_4 & w_5 \\ \hline t_1 & 0 & 1 & 3 & 2 & 3 \\ t_2 & 0 & 2 & 4 & 1 & 3 \\ t_3 & 0 & 2 & 4 & 1 & 4 \\ t_4 & 0 & 4 & 4 & 3 & 1 \end{array}$$

$$s = \operatorname*{argmax}_{i} c_{i,K} = 1$$

<s>  w1   w2   w3   w4   w5

# Backward pass

$$D = \begin{array}{c|c|c|c|c|c|} & w_1 & w_2 & w_3 & w_4 & w_5 \\ \hline t_1 & 0 & 1 & 3 & 2 & \boxed{3} \\ \hline t_2 & 0 & 2 & 4 & 1 & 3 \\ \hline t_3 & 0 & 2 & 4 & \boxed{1} & 4 \\ \hline t_4 & 0 & 4 & 4 & 3 & 1 \\ \hline \end{array}$$

```
<s>  w1   w2   w3   w4   w5
                              t1
```

```
<s>  w1   w2   w3   w4   w5
                         t3 ← t1
```

# Backward pass

$$D = \begin{array}{c|c|c|c|c|c|} & \text{w}_1 & \text{w}_2 & \text{w}_3 & \text{w}_4 & \text{w}_5 \\ \hline \text{t}_1 & 0 & 1 & 3 & 2 & 3 \\ \hline \text{t}_2 & 0 & 2 & 4 & 1 & 3 \\ \hline \text{t}_3 & 0 & 2 & 4 & 1 & 4 \\ \hline \text{t}_4 & 0 & 4 & 4 & 3 & 1 \\ \hline \end{array}$$

\<s\>   w1   w2   w3   w4   w5

$t_1 \leftarrow t_3 \leftarrow t_1$

# Backward pass

$$D = \begin{array}{c|ccccc} & w_1 & w_2 & w_3 & w_4 & w_5 \\ \hline t_1 & 0 & 1 & 3 & 2 & 3 \\ t_2 & 0 & 2 & 4 & 1 & 3 \\ t_3 & 0 & 2 & 4 & 1 & 4 \\ t_4 & 0 & 4 & 4 & 3 & 1 \end{array}$$

&lt;s&gt;   w1   w2   w3   w4   w5

$t_1 \leftarrow t_3 \leftarrow t_1$

# Backward pass

$$D = \begin{array}{c|ccccc} & w_1 & w_2 & w_3 & w_4 & w_5 \\ \hline t_1 & 0 & 1 & 3 & 2 & 3 \\ t_2 & 0 & 2 & 4 & 1 & 3 \\ t_3 & 0 & 2 & 4 & 1 & 4 \\ t_4 & 0 & 4 & 4 & 3 & 1 \end{array}$$

<s>   w1   w2   w3   w4   w5

$t_3 \leftarrow t_1 \leftarrow t_3 \leftarrow t_1$

# Backward pass

$$D = \begin{array}{c|ccccc} & w_1 & w_2 & w_3 & w_4 & w_5 \\ \hline t_1 & 0 & 1 & 3 & 2 & 3 \\ t_2 & 0 & 2 & 4 & 1 & 3 \\ t_3 & 0 & 2 & 4 & 1 & 4 \\ t_4 & 0 & 4 & 4 & 3 & 1 \end{array}$$

&lt;s&gt;  w1   w2   w3   w4   w5

$t_3 \leftarrow t_1 \leftarrow t_3 \leftarrow t_1$

# Backward pass

$$D = $$

| | $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ |
|---|---|---|---|---|---|
| $t_1$ | 0 | 1 | 3 | 2 | 3 |
| $t_2$ | 0 | 2 | 4 | 1 | 3 |
| $t_3$ | 0 | 2 | 4 | 1 | 4 |
| $t_4$ | 0 | 4 | 4 | 3 | 1 |

&lt;s&gt;  w1   w2   w3   w4   w5

$\pi$  $\leftarrow t_2 \leftarrow t_3 \leftarrow t_1 \leftarrow t_3 \leftarrow t_1$

# Named Entity Recognition

# Named Entities

- In its core usage, a named entity means anything that can be referred to with a proper name

- Most common 4 tags:
  - PER (Person): "Marie Curie"
  - LOC (Location): "New York City"
  - ORG (Organization): "Stanford University"
  - GPE (Geo-Political Entity): "Boulder, Colorado"
  - Often multi-word phrases
    - But the term is also extended to things that aren't entities: dates, times, prices

- The task of named entity recognition (NER)
  - find spans of text that constitute proper names
  - tag the type of the entity

# NER output

Citing high fuel prices, [ORG **United Airlines**] said [TIME **Friday**] it has increased fares by [MONEY **$6**] per round trip on flights to some cities also served by lower-cost carriers. [ORG **American Airlines**], a unit of [ORG **AMR Corp.**], immediately matched the move, spokesman [PER **Tim Wagner**] said. [ORG **United**], a unit of [ORG **UAL Corp.**], said the increase took effect [TIME **Thursday**] and applies to most routes where it competes against discount carriers, such as [LOC **Chicago**] to [LOC **Dallas**] and [LOC **Denver**] to [LOC **San Francisco**].

# Why NER?

- Sentiment analysis
  - consumer's sentiment toward a particular company or person?
- Question Answering
  - answer questions about an entity?
- Information Extraction
  - Extracting facts about entities from text

# Why NER is hard

- Segmentation
  - In POS tagging, no segmentation problem since each word gets one tag
  - In NER we must find and segment the entities!

- Type ambiguity

[PER Washington] was born into slavery on the farm of James Burroughs.
[ORG Washington] went up 2 games to 1 in the four-game series.
Blair arrived in [LOC Washington] for what may well be his last state visit.
In June, [GPE Washington] passed a primary seatbelt law.

# BIO Tagging

- How can we turn this structured problem into a sequence problem like POS tagging, with one label per word?

- [PER Jane Villanueva] of [ORG United], a unit of [ORG United Airlines Holding], said the fare applies to the [LOC Chicago ] route

# BIO Tagging

- [PER Jane Villanueva] of [ORG United], a unit of [ORG United Airlines Holding], said the fare applies to the [LOC Chicago ] route

| Words | BIO Label |
|---|---|
| Jane | B-PER |
| Villanueva | I-PER |
| of | O |
| United | B-ORG |
| Airlines | I-ORG |
| Holding | I-ORG |
| discussed | O |
| the | O |
| Chicago | B-LOC |
| route | O |
| . | O |

Now we have one tag per token!!!

# BIO Tagging

- B: token that *begins* a span

- I: tokens *inside* a span

- O: tokens outside of any span

- \# of tags (where n is #entity types):
  - 1 O tag,
  - *n* B tags,
  - *n* I tags

- total of *2n+1*

| Words | BIO Label |
|---|---|
| Jane | B-PER |
| Villanueva | I-PER |
| of | O |
| United | B-ORG |
| Airlines | I-ORG |
| Holding | I-ORG |
| discussed | O |
| the | O |
| Chicago | B-LOC |
| route | O |
| . | O |

# BIO Tagging variants: IO and BIOES

- [PER Jane Villanueva] of [ORG United], a unit of [ORG United Airlines Holding], said the fare applies to the [LOC Chicago ] route

| Words | IO Label | BIO Label | BIOES Label |
|---|---|---|---|
| Jane | I-PER | B-PER | B-PER |
| Villanueva | I-PER | I-PER | E-PER |
| of | O | O | O |
| United | I-ORG | B-ORG | B-ORG |
| Airlines | I-ORG | I-ORG | I-ORG |
| Holding | I-ORG | I-ORG | E-ORG |
| discussed | O | O | O |
| the | O | O | O |
| Chicago | I-LOC | B-LOC | S-LOC |
| route | O | O | O |
| . | O | O | O |

# Standard algorithms for NER

- Supervised Machine Learning given a human-labeled training set of text annotated with tags
  - Hidden Markov Models
  - Conditional Random Fields (CRF)/ Maximum Entropy Markov Models (MEMM)
  - Neural sequence models (RNNs or Transformers)
  - Large Language Models (like BERT), finetuned